

Thursday, March 19, 2020 6:44 PM

- Smoothing parameter,  $h$ , free method
- Reduction of Computational burden
- Conditions may not be applicable for many applications

$$f(x) = \sum_{i=1}^c p_i \cdot g(x; \theta_i), \quad c \ll n$$

weight      component density with parameter  $\theta_i$

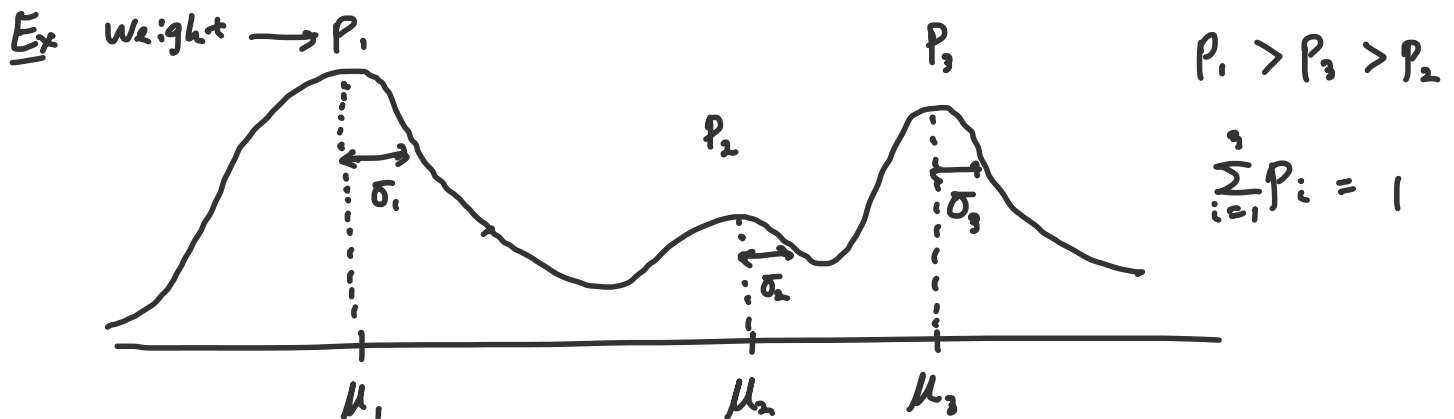
$\sum_{i=1}^c p_i = 1, p_i > 0$  at  $x$

- The component density can be either continuous or discrete.
- # of component density,  $c$ , need to be known in advance  
usually normal density
- $p_i$  and  $\underline{\theta}_i$ ,  $i=1, \dots, c$  are to be estimated.

$$\hat{f}_{\text{FM}}(x) = \sum_{i=1}^c \hat{p}_i \cdot \underbrace{\phi(x; \hat{\mu}_i, \hat{\sigma}_i^2)}_{\text{normal density function}}$$

- The observation  $x$  comes from one of these normal distribution but we do not observe to which component it belongs.
- Parameters to be estimated ;  $p_1, p_2, \dots, p_{c-1}, \mu_1, \mu_2, \dots, \mu_c$   
 $\sigma_1^2, \sigma_2^2, \dots, \sigma_c^2$

Total # of parameters =  $3c - 1$



Example 9.8 Univariate Finite Mixture model with 3 normal Components

$$f(x) = \underset{p_1}{0.3} \times \underset{\mu_1}{\phi}(x; \underset{\sigma_1}{-3}, \underset{p_2}{1}) + \underset{p_2}{0.3} \times \underset{\mu_2}{\phi}(x; \underset{\sigma_2}{0}, \underset{p_3}{1}) + \underset{p_3}{0.4} \times \underset{\mu_3}{\phi}(x; \underset{\sigma_3}{2}, \underset{\sigma_3}{0.5})$$

See the sample code.

- Once a sample data set is given, we need to identify the number of components,  $c$ , and the Component density.
- EM (Expectation Maximization) algorithm can be used

to estimate  $\mu$ 's,  $\sigma$ 's and  $\pi$ 's.

## Finite Mixture vs Kernel density estimation

$$\hat{f}_{FM}(x) = \sum_{i=1}^C \hat{\pi}_i \cdot g(x; \hat{\theta}_i) \quad \text{vs} \quad \hat{f}_{ker}(x) = \frac{1}{n} \sum_{i=1}^n K_h(x - x_i)$$

$\nwarrow$  weight       $\downarrow$  component density       $\nwarrow$  equal weight       $\downarrow$  kernel density as a function of  $h$

• Kernel est is a special case of FM where  $C = n$

## Parameter Estimation

EM algorithm (Dempster, Laird and Rubin '77)

$\downarrow$  Expectation       $\searrow$  Maximization

Optimizing likelihood function with missing data

[csfminix]

## Iterative Procedure (Univariate)

Normal component density case

Step 1: Determine  $C$  (exploratory data analysis can be

used)

Step 2 : Determine initial guess for  $\hat{p}_i$ 's,  $\hat{\mu}_i$ 's and  $\hat{\sigma}_i^2$ 's

Step 3 : Calculate the estimated posterior prob. that  $x_j$   
 $j=1, \dots, n$  belongs to ↓  
data

$$\hat{z}_{ij} = \frac{\hat{p}_i \cdot \phi(x_j; \hat{\mu}_i, \hat{\sigma}_i^2)}{\hat{f}(x_j)}$$

↘  $\sum_{k=1}^c \hat{p}_k \cdot \phi(x_j; \hat{\mu}_k, \hat{\sigma}_k^2)$

Prob. that  $x_j$  belongs to  $\phi(\hat{\mu}_i, \hat{\sigma}_i^2)$ ,  $i^{\text{th}}$  component.

Step 4 : Update

$$\hat{p}_i = \frac{1}{n} \sum_{j=1}^n \hat{z}_{ij}$$

$$\hat{\mu}_i = \frac{1}{n} \sum_{j=1}^n \frac{\hat{z}_{ij} x_j}{\hat{p}_i} \quad : \text{weighted average}$$

$$\hat{\sigma}_i^2 = \frac{1}{n} \sum_{j=1}^n \frac{\hat{z}_{ij} (x_j - \hat{\mu}_i)^2}{\hat{p}_i}$$

Step 5 : Repeat step 3 and 4 until converge

< See Example 9.13 >

Generating r.v. from the finite Mixture model

Step 1 : Get the FM model :  $p_i, g_i(x; \theta_i), i=1, \dots, c$

Step 2 : Generate  $n$  r.v. from  $U(0, 1)$  ;  $u_1, u_2, \dots, u_n$

Step 3 : Count # of  $u$ 's in  $[0, p_1) \Rightarrow n_1$   
 " " in  $[p_1, p_1+p_2) \Rightarrow n_2$   
 $\vdots$   
 Count # of  $u$ 's in  $[p_1+\dots+p_{c-1}, 1] \Rightarrow n_c$

Step 4 : Generate  $n_1$  r.v. ,  $x_1, \dots, x_{n_1}$  from  $g_1$   
 $n_2$  r.v. ,  $x_{n_1+1}, \dots, x_{n_1+n_2}$  "  $g_2$   
 $\vdots$   
 $n_c$  " " "  $g_c$

< See example 9.13 >

In-class assignment

normal

- Write a procedure to generate r.v. from Kernel density estimation (similar to the procedure for FM)
- (a) Create artificial 3-term mixture data with  $n = 1500$  of which  
 $n_1 = 200$  from  $N(5, 3)$

$n_2 = 800$  from  $N(10, 1.5^2)$

$n_3 = 500$  from  $N(15, 2^2)$

This is your raw data.

(b) Find the finite mixture model using EM  
[csfinmix]

Use the initial  $p = [\frac{1}{3}, \frac{1}{3}, \frac{1}{3}]$

$\mu = [4, 12, 16]$

$\sigma = [2, 1, 2]$

(c) Now, generate  $n = 1500$  r.v. from the finite mixture model.

(d) Draw density histogram of the raw data in (a) and r.s. in (c). Compare.

3. Using the same raw data in #2, estimate the density using normal kernel density estimation method. Superimpose the density estimation on the histogram of the raw data.

Repeat (c) and (d) in #2.